



Lunt, G., Day, MA., & Wilson, RE. (2006). *Enhancing motorway traffic data with novel vehicle re-identification algorithms*.
<http://hdl.handle.net/1983/497>

Early version, also known as pre-print

[Link to publication record in Explore Bristol Research](#)
PDF-document

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

Enhancing motorway traffic data with novel vehicle re-identification algorithms

ITS World Congress Scientific Paper Submission
January 2006

George Lunt, TRL Limited, glunt@trl.co.uk
Mark Day, University of Bristol, md2409@bristol.ac.uk
R. Eddie Wilson, University of Bristol, RE.Wilson@bristol.ac.uk

1. Introduction

This paper describes recent work in vehicle re-identification procedures applied to individual vehicle data (IVD) from highway inductance loop systems. Briefly, vehicle re-identification is the process of matching data from sensor systems at different locations so as to track an individual vehicle down the highway. Perhaps the most accurate method of re-identification is automatic number plate recognition (ANPR) [1]. However, the idea in this paper, pioneered by Coifman and collaborators [2],[3] is that re-identification may also be accomplished more cheaply by the better analysis of data from existing inductance loop systems. In the past, the accuracy and proportion of vehicles that could be re-identified in this way has been quite limited. The advances that we report here, which are due both to improved inductance loop systems and to more refined data analysis techniques, push the re-identification rate close to 100%. Papers [4] and [5] discuss early designs and progress in this research programme.

The work described here has been carried out by TRL in collaboration with the University of Bristol, and is based on data from the United Kingdom's Motorway Incident Detection and Automatic Signalling (MIDAS) system [6], provided by kind permission of the UK Highways Agency (HA).

The MIDAS system protects slow moving traffic by monitoring traffic conditions, and by setting signals (such as temporary speed limits) according to pre-defined rules. Its detection system consists of inductance loop pairs which are connected to outstations that contain a smart signal processing and communications box. This set-up uses the magnetic field of each passing vehicle to capture the Individual Vehicle Data (IVD) listed in Table 1. However, in normal operation, outstations do not store IVD, but rather combine measurements from multiple vehicles into one minute averages, which are sent to a central control computer for processing. However, it is possible to intercept the IVD before it is discarded, by connecting a laptop computer inside the roadside cabinet of each outstation of interest.

The inductance loop sites operated under MIDAS are usually installed at 500 metre intervals along the highway, which is a comparable spacing to that used in the Berkeley Highway Laboratory's well-known I-880 data set [7], used by [2],[3]. However, on the Active Traffic Management (ATM) section of the M42 motorway near the city of Birmingham UK, the interval is reduced to 100 metres so as to help improve incident response times [8]. Between 8th and 10th December 2003, TRL collected MIDAS IVD at six consecutive northbound sites on this three-lane stretch of motorway, between the exit and entry slip of Junction 6. During the 48 hour exercise, approximately 90000 vehicles drove through this section of road and their IVD was used to develop and validate the re-identification algorithms discussed later in the paper.

Section 2 describes a package which was developed to visualise the data from this project. It becomes apparent that the data is of such fine resolution that a human operator may re-identify 'by eye' almost all of the vehicle records across all six sites, with doubt only arising in some rare situations where the traffic is heavy and a

number of lane changes occur simultaneously. Therefore a small *hand-matched* data set was built for use as the *ground truth* in the evaluation of the automatic re-identification algorithms. As a cross-check, a micro-simulation package was also used to generate artificial IVD, for which (by construction) the ground truth re-identification was known.

Measurement	Units
Unique ID	-
Time-stamp	Seconds
Speed	Kilometres per hour
Length	Centimetre
Front-to-front headway	Tenths of a second
Lane	Integer lane number

Table 1. IVD measurements. All data has integer type

Section 3 discusses automatic re-identification algorithms. Since the full details are proprietary, only the general procedure is outlined. In contrast to Coifman [2],[3] who uses vehicle length as the chief piece of data to match, the algorithms discussed in this paper use velocities, lengths and time-stamps. Further, here the interval between loop pair sites is so small that there are typically few lane changes between them, and consequently headway sequences change very little. Both these properties are exploited in the algorithms here, which thus improve significantly on the 60% re-identification rate of [2],[3].

Finally, note that if one can re-identify a vehicle through a whole sequence of loop pair sites, then one has in effect re-constructed its trajectory. If this process can be carried out for a large proportion of the vehicles on the highway, then the resulting data would hold a wealth of information on inter-vehicle dynamics, and could help to improve our understanding of the mechanisms behind flow breakdown. Furthermore, if IVD could be obtained and re-identified in real time, then improved incident detection techniques algorithms would result. A summary of future applications is given in Section 4.

2. Visualisation and Manual Re-identification

To help with the development of the re-identification algorithms, a Matlab [9] interface has been developed to display the IVD, see Figure 1. This figure displays a 20 second scope plot, in which the six M42 sites are labelled A,B,C,D,E, and F in downstream order, so that vehicles drive past site A first and past site F last. Each site is represented by a new axis, or sub-plot, on which the IVD from that site is displayed.

The vertical axis on each sub-plot is time, and vehicle records are depicted by filled rectangles, labelled with speeds in kph, whose sizes are representative of lengths and whose positions are determined by lane numbers and time-stamps. In practice, since time-stamps are only accurate to one second, headway information (accurate to tenths of a second) is used to refine the relative positions in which vehicles' rectangles are plotted.

In Figure 1, the vertical axis has been offset by 4 seconds from site to site so that a vehicle travelling at 90 kph should maintain the same vertical level from sub-plot to sub-plot, provided the site spacing is precisely 100 metres (in fact, this is only

approximately so). The labelling of the vertical time axis refers to the left-most sub-plot.

The scope plot can be thought of as rather like a set of overhead views of the highway at each of the six sites, in which vehicles are driving up the page but also progressing from left to right (downstream) in the sequence of sub-plots. Note that these plots depict the UK situation where slow traffic (lane 1) drives on the left. Of course, the overhead view analogy uses the vertical time axis as a space-like coordinate and consequently breaks down if the vehicles have widely varying velocities.

It should now be clear how the pattern of vehicle records is replicated from site to site, and Figure 2 adds manual re-identification information to Figure 1. The vehicles bounded by the horizontal lines in the first sub-plot have been manually re-identified across the subsequent five sites, and the vehicles boxed by rectangles are thought not to have been involved in lane changes and so their order is preserved.

In contrast, it seems there must be a lane change at the rear of the boxed groups, and there appears to be only one plausible way in which the re-identification can be performed. The circled vehicle is recorded in lane 2 at site A, lane 3 at sites B-E, and then again in lane 2 at site F. This particular vehicle has moved out of lane 2 in order to overtake a slower moving platoon, but has then pulled back in when a much faster vehicle approaches from behind in lane 3. This example illustrates the fine level of microscopic detail that can be extracted from re-identified IVD. In particular, it seems plausible that a large re-identified data set could be used to gather new microscopic behavioural statistics such as gap acceptance distributions and so help to improve the state of the art in simulation software.

The Matlab interface has been further developed so that the user may perform manual re-identification and store results using standard point-and-click operations on the scope plot. This procedure has been used to build a number of 'hand-matched' sets within the IVD, containing about 3000 vehicle trajectories in total. Examination of this manually re-identified data has helped translate the intuitive techniques that human operators use into a formal set of rules for automatic re-identification algorithms. Furthermore, the hand-matched set is used as ground truth for testing the efficacy of those algorithms. Work is presently under way to hand-match more data so that the training of algorithms and their evaluation may be performed independently.

In order to validate the accuracy of the manual re-identification process itself, it was necessary to have an IVD set for which the ground truth was independently available. The approach undertaken involved using SISTM [10] (a micro simulation model whose name stands for SImulation of Strategies for Traffic on Motorways) to create artificial IVD with the same six-site structure as the real data. A unique reference number for each vehicle is appended to the artificial IVD during simulation and matching these references across sites gives the ground truth re-identification, which is hidden from the operator during the hand-matching exercise.

Of the 1320 SISTM vehicles analysed in this exercise, 21 were hand-matched incorrectly, which corresponds to an error rate of 1.6%. This figure is rather high, since the accuracy of the automatic re-identification algorithms (i.e. their discrepancy with hand-matched data) is itself of the order of 1% (see Section 3). However, it appears that hand-matching real IVD is easier than hand-matching SISTM data, because in SISTM, the velocity fluctuations of individual vehicles over 100 metres are larger than in real data. Consequently, SISTM's headway patterns do not replicate so closely from site to site and there are even instances where a vehicle leaves a platoon by changing lane and the resulting gap is closed before the next site. This

kind of behaviour does not occur in real traffic. As a result it is thought that the error rate for hand-matching real IVD is much lower than the 1.6% rate for SISTM data.

This discussion highlights the difficulties in developing accurate microscopic simulation models when there is not sufficient empirical data to allow for robust calibration and validation at the microscopic level. An interesting future exercise would concern the comparison of the dynamics of a range of microscopic simulation models against re-identified IVD.

3. Automatic Re-identification Algorithms

This section describes two algorithms that have been developed to automatically re-identify IVD. MIDAS sites spaced at 100 metre intervals are unique to the ATM section of the M42 motorway, and hence it was considered desirable to develop not only an algorithm for 100 metre data, but also a separate algorithm that deals with the standard MIDAS loop spacing of 500 metres used throughout the rest of the UK. Since IVD has been captured from six consecutive sites on the M42, 500 metre algorithms may be tested by using the data from the first and last site only. This procedure is much better than testing against 500 metre data captured elsewhere, because the data from intermediate loops may be used during the hand-matching process to construct an almost perfect ground truth set. In practice, it is very difficult to hand-match 500 metre data if no intermediate information is available, because vehicles may re-order significantly over this distance.

The basic design for both algorithms is very similar and is described briefly below. It should be noted that precise details such as optimal tolerances etc. are proprietary and hence have been omitted from the discussion.

The algorithms start by taking each vehicle from an upstream site and they consider downstream vehicles as possible matches if they meet certain criteria (e.g. within a certain speed and length tolerance and recorded within a required time window). These criteria are strict for the first 'sweep'. Vehicles that are uniquely matched to each other (i.e. no other vehicles consider them as a possible match) are flagged as definite matches.

Following this initial sweep, vehicles with several possible matches are re-considered. For these vehicles, the upstream and downstream headways to the closest previously matched vehicles are calculated. Possible matches at the downstream loop, and their associated front and back headways with the same matched vehicle, are compared. If comparison is within a certain tolerance (which for the first execution is tight), then the pair is marked as a possible match. After considering all vehicles again, unique matches are flagged as definite matches. This process is repeated a number of times, re-introducing unmatched vehicles and relaxing the headway and other tolerances at each repetition. The 100 metre algorithm has stricter criteria than the 500 metre algorithm at each iteration since the smaller loop spacing allows for more accurate predictions in the arrival times, velocities and headways at the downstream site, and there is also less lane changing between consecutive sites.

Note that initially the algorithms seek matches with no lane changes, but this requirement is relaxed in the latter part of the algorithm and generally speaking lane changes are correctly identified except in some cases where several of them occur simultaneously.

The 100 metre algorithm performs an additional final step by considering all remaining unmatched vehicles. For all possible downstream matches, the algorithm applies a matching score based on the length difference, velocity difference,

expected arrival time, and the number of lane changes made by the vehicle. These scores are stored in a *match matrix*.

A hypothetical match matrix for upstream vehicles A1,B1,C1,D1,E1 on downstream vehicles A2,B2,C2,D2,E2 could be as shown in Table 2. The numbers within the match-matrix are the scores, which represent the likelihood for each putative match, with high scores representing good matches, and blank entries representing no possible match.

		downstream vehicles				
		A2	B2	C2	D2	E2
upstream vehicles	A1	8	6	3		
	B1	7	8	2	1	
	C1		7	7		
	D1		8	9	5	7
	E1			9	9	9

Table 2. Hypothetical match matrix (high numbers represent a good match; low numbers represent a bad match)

The combination of matches is then found such that the total sum of the scores is maximised. This is known as the *weighted bipartite matching problem* and is solved by the Hungarian Algorithm [11]. In this particular example the solution is (A1,A2), (B1,B2), (C1,C2), (D1,E2), (E1,D2) with a maximised total score of 39.

As discussed in Section 2, a variety of hand-matched data (regarded as the ground truth) is used to evaluate the performance of each of the two algorithms. Two basic error statistics have been used, namely the *re-identification rate* and the *failure rate*. For a set of n vehicles considered for matching where a are matched correctly and b are matched incorrectly, the error statistics are defined as follows:

$$\begin{aligned} \text{re-identification rate} &= \frac{a+b}{n}, \\ \text{failure rate} &= \frac{b}{a+b}. \end{aligned}$$

The re-identification rate is thus a measure of the proportion of vehicles matched (correctly or incorrectly) by the algorithm and the failure rate measures the proportion of those matches which are incorrect. As yet no attempt has been made to derive an overall performance statistic (e.g. define an objective function that combines the re-identification and failure rates) since the relative importance of unmatched vehicles and of false positives will vary according to the application.

Results are summarised in Table 3. Since the 100 metre algorithm invokes the final Hungarian Algorithm step, the re-identification rate is total (100%). In contrast, the 500 metre algorithm does not perform this final step and so leaves some vehicles unmatched (re-identification rate of 94%). The failure rate of the 100 metre algorithm is very low, at order 1%. This figure cannot be quoted any more accurately due to the relatively small size of the hand-matched sets and the possibility (discussed earlier) that there are errors in the hand-matched data itself.

Not surprisingly, the 500 metre algorithm has a much higher failure rate (circa 7%) than the 100 metre algorithm. Note that by applying the 100 metre algorithm five times in succession, and assuming independence of its failures, one would achieve a re-identification rate of 100% and a failure rate of only 5% approximately, i.e., one would out-perform the 500 metre algorithm. This statistic underlines the value of having loop sites close together when trying to follow individual vehicles a long distance down the highway.

Algorithm name	Re-identification rate	Failure rate
(1) 100 metre Algorithm	100%	1%
(2) 500 metre Algorithm	94%	7%

Table 3. Performance of Automatic Re-Identification Algorithms

4. Future Applications and Conclusions

This paper has described the design and performance of two vehicle re-identification algorithms applied to MIDAS IVD inductance loop data. Some graphical visualisation software has also been developed and the performance of the algorithms has been evaluated using a manually re-identified data set.

The 100 metre algorithm in particular re-identifies all vehicles in the test set with an error rate of around only 1%, and is ready for use in applied research projects (potential applications listed below). In contrast, the 500 metre algorithm does not re-identify all vehicles and moreover has a false positive rate of around 7%, but nevertheless 85%-90% of all vehicles are correctly re-identified. This is a substantial improvement on the 60% rate reported in [2],[3] whose sites have a similar spacing. This improvement is probably due to the high quality of the double loops employed in the UK and the variety of fields (time-stamp, velocity, length and headway) that have been used in the algorithm.

Further work is continuing in improving both algorithms and also testing them against more extensive IVD sets. The sites used here, being mid-junction (after the off-slip and before the on-slip) are not generic, and moreover the flow levels recorded never exceeded more than about 4000 veh/hr across the three lanes of the carriageway. As a consequence, the traffic flow was freely flowing and it will be necessary to validate the algorithms in more turbulent conditions (e.g. on a main link, close to a junction and operating close to capacity). However, this validation requires more extensive data collection.

Some potential future applications for re-identified data are listed below.

1. Improvements in Microscopic Model Calibration

There is presently no comprehensive empirical data set for *microscopic* highway traffic dynamics. Hence modellers are forced to guess their own relationships concerning vehicle acceleration and lane changing behaviour, and validation is only possible at the level of the emergent macroscopic dynamics. A large re-identified IVD set will go a long way towards rectifying this problem.

2. Incident Detection

If IVD were available in real-time, then real-time re-identification could be used to power incident detection algorithms, capable of detecting the fast forward travelling shock-wave (lack of vehicles) rather than the resulting queue which propagates upstream rather slowly, see [12].

3. Improved Estimates of Link Capacities

Lane changing is widely believed to be one of the main influences on link capacity. However, link capacity formulae (as for example, in COBA) do not take into account lane changing rates in their calculations. This is because lane-changing rates have

not been collected in the required quantity to derive statistically robust relationships. However, lane changing rates could be derived from re-identified IVD, and such capacity relationships could consequently be reviewed and improved.

4. Improved Understanding of Flow Breakdown

It is widely believed that when motorway sections are operating close to capacity, perturbations in the traffic flow are the principal cause of flow breakdown and shock-wave generation. However, the exact (microscopic) conditions leading to flow-breakdown are not known. Re-identified IVD will help improve knowledge in this area, and make it easier to provide more targeted solutions to reduce the causes of congestion.

5. Loop Tuning

Loop sites that give very different length measurements for the same re-identified vehicle are poorly tuned. The re-identified IVD could be used to help diagnose which sites were poorly tuned, and improve the fidelity of the traffic measurements coming from them. This application would be particularly useful if IVD were available in real-time as drift could also be identified and corrected in real time.

To conclude, the output from the algorithms described here provides traffic data at a level of detail and quantity never previously possible. This provides the framework for more detailed research into highway traffic behaviour and the specific causes of congestion. The data also provides the opportunity for incorporation into future intelligent transport systems, e.g. powering improved incident detection algorithms and automatic loop calibration.

References

- [1] Frith B, Pearce D (2002). *Driver Information on Journey Time Variability Generated using ANPR Data*. 2002 RTIC conference proceedings.
- [2] Coifman B, Cassidy M (2002). *Vehicle Reidentification and Travel Time Measurement on Congested Freeways*. Transportation Research: Part A, 2002, vol 36, no 10, 2002, pp 899-917.
- [3] Coifman, B (1998). *Vehicle Reidentification and Travel Time Measurement in Real-Time on Freeways Using the Existing Loop Detector Infrastructure*. Transportation Research Board, 1998, pp 181-191.
- [4] Lunt G.M., Wilson R.E. (2003). *New Data Sets and Improved Models of Highway Traffic*. UTSG 35th Annual Conference, Volume 2
- [5] Lunt G.M. *Vehicle Re-Identification Using Induction Loop Data*. ECTRI-FEHL-FERSI Young Researchers Seminar 2005.
- [6] MIDAS Information from the UK Highways Agency's website
<http://www.highways.gov.uk/news/pressrelease.aspx?pressreleaseid=686>
- [7] Skabardonis, A, et al (1996). *I-880 Field Experiment: Data-Base Development and Incident Delay Estimation Procedures*. Transportation Research Record 1554, TRB, 1996, pp 204-212.
- [8] Highways Agency. 2005. *Active Traffic Management (ATM) Project M42 Junctions 3A-7*. <http://www.highways.gov.uk/knowledge/1361.aspx>
- [9] The Mathworks website, developers of the Matlab software environment.
<http://www.mathworks.co.uk/>
- [10] SISTM (1993). *A Motorway Simulation Model*. Leaflet LF2061. TRL Limited.
- [11] Galil Z (1986). *Efficient Algorithms for finding maximum matching in graphs*. ACM Computing Surveys (CSUR) Volume 18, Issue 1
- [12] Coifman B (2003). *Identifying the Onset of Congestion Rapidly with Existing Traffic Detectors*. Transportation Research:Part A, vol 37, no 3, 2003, pp 277-291.

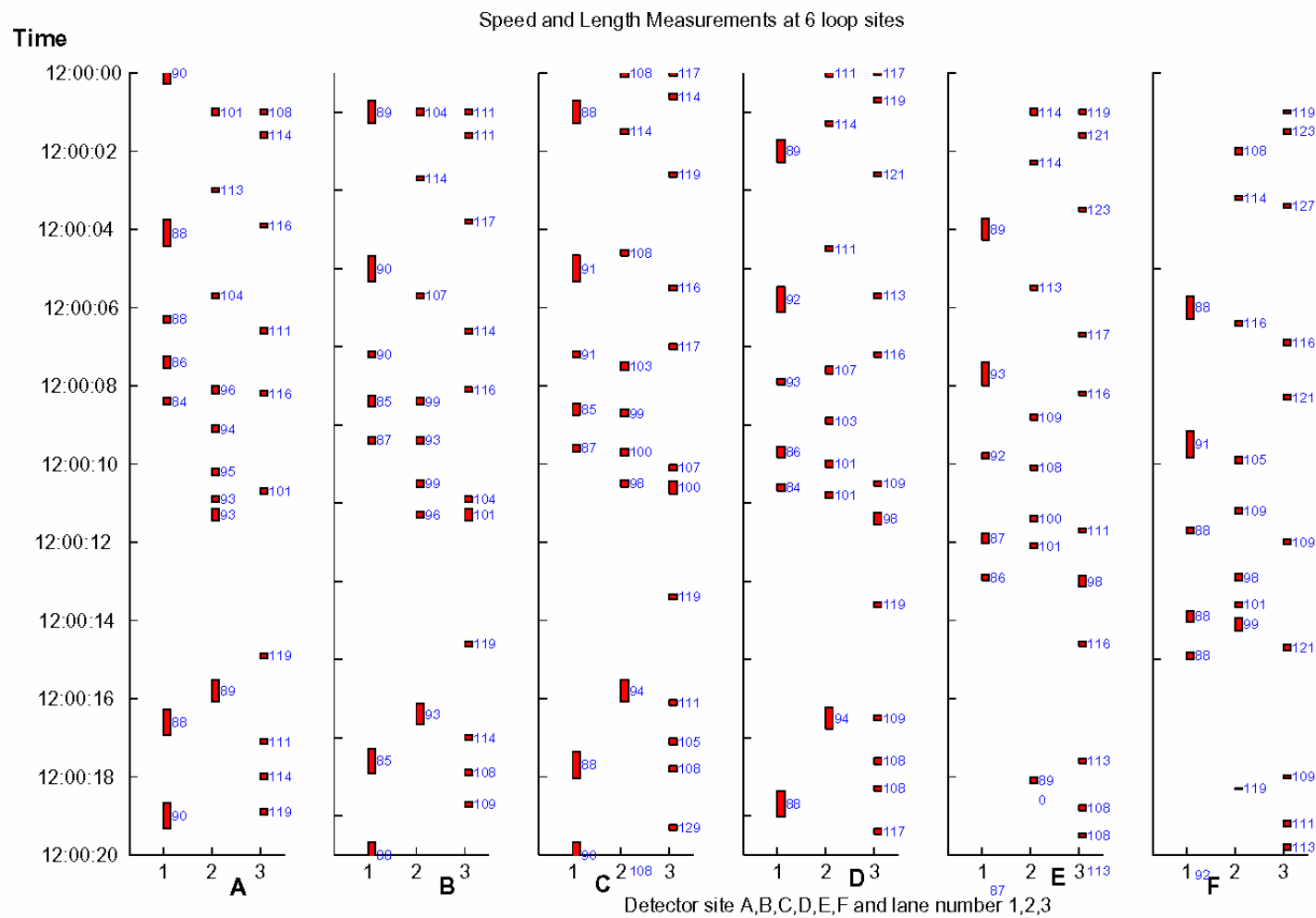


Figure 1. Vehicle records at six consecutive loop sites, 100 metres apart. M42 Junction 6.

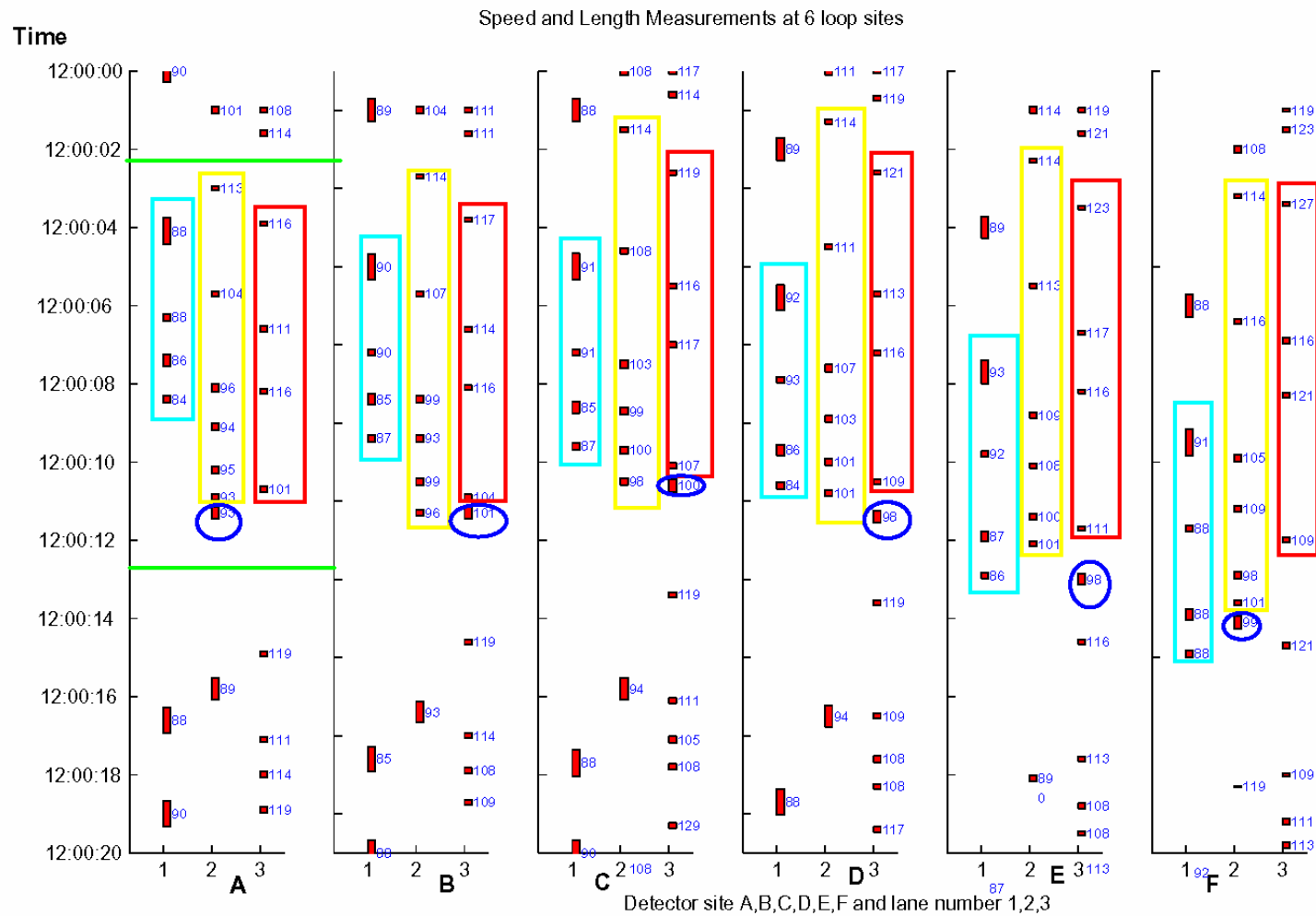


Figure 2. Re-identified vehicle records at six consecutive loop sites, 100 metres apart. M42 Junction 6.